

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
28 November 2002 (28.11.2002)

PCT

(10) International Publication Number
WO 02/095730 A1

(51) International Patent Classification⁷: **G10L 15/20**

(21) International Application Number: **PCT/GB02/02197**

(22) International Filing Date: **20 May 2002 (20.05.2002)**

(25) Filing Language: **English**

(26) Publication Language: **English**

(30) Priority Data:
0112319.9 **21 May 2001 (21.05.2001)** **GB**

(71) Applicant (*for all designated States except US*):
QUEEN'S UNIVERSITY OF BELFAST [GB/GB];
8 Malone Road, Belfast BT9 5BN (GB).

(72) Inventor; and

(75) Inventor/Applicant (*for US only*): **MING, Ji [GB/GB];**
10 Rosssdale Heights, Belfast BT8 6XZ (GB).

(74) Agent: **MURGITROYD & COMPANY;** Scotland
House, 165-169 Scotland Street, Glasgow G5 8PL (GB).

(81) Designated States (*national*): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NO, NZ, OM, PH, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VN, YU, ZA, ZM, ZW.

(84) Designated States (*regional*): ARIPO patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE, TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

Published:

— *with international search report*

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

(54) Title: **INTERPRETATION OF FEATURES FOR SIGNAL PROCESSING AND PATTERN RECOGNITION**

(57) Abstract: A method of interpretation of features for signal processing and pattern recognition provides a model in which the pattern or signal to be interpreted is considered as a set of N observations, M of which are corrupt, and a disjunction is performed over all possible combinations of N different values (1,...,N) taken N-M at a time. The value of M defines the order of the model, and is determined using an optimality criterion which chooses the order that corresponds to a clean signal based on comparing the state duration probability of the signal or pattern to be interpreted with that of a clean signal.

WO 02/095730 A1

1 Interpretation of Features for Signal Processing and
2 Pattern Recognition

3
4 The present invention relates to interpretation of
5 features for signal processing and pattern
6 recognition, and particularly to speech recognition
7 subjected to partial, unknown frequency-based
8 corruption.

9
10 Partial frequency-band corruption may account for
11 the effect of a family of real-world noises, for
12 example, a telephone ring, a car horn, a siren or a
13 random channel tone, which usually have a band-
14 selective characteristic and thus affect only
15 certain parts of the speech frequency band. There
16 may be two different ways to deal with this type of
17 noise corruption for robust speech recognition.
18 Firstly, we may use the conventional noise filtering
19 or feature/model compensation techniques to remove
20 the noise component from the input signal, or to
21 adapt the model to the noisy environment. Each of
22 these techniques assumes the availability of certain

CONFIRMATION COPY

1 knowledge of the noise or environment. The required
2 knowledge may include, for example, the spectral or
3 cepstral characteristics of the noise for noise
4 filtering or feature selection, a stochastic model
5 of the noise for noise compensation and an extra set
6 of training data in the new environment for model
7 adaptation.

8
9 The second possible way of dealing with this partial
10 corruption is to base the recognition mainly on
11 information from the clean frequency bands, by
12 throwing away the noisy bands, or by making these
13 bands play a less significant role in recognition,
14 i.e., the missing feature method.

15
16 This recognition is made possible due to redundancy
17 in the spectral characteristics of speech. This
18 method is of interest because there can be
19 situations where removing the noise from the input
20 signal may prove difficult, due to the lack of
21 sufficient knowledge about the noise. This lack of
22 knowledge may be experienced, for example, when an
23 unknown unexpected noise occurs in the middle of
24 utterance. A better system may be a combination of
25 these two methods, i.e., using the noise reduction
26 technique to remove the noise with a known or
27 stationary characteristic, and exploiting the
28 redundancy in the speech signal to get around the
29 noise with an unknown or time-varying nature. The
30 present invention focuses on the second method, but
31 we use a simple example to demonstrate the advantage
32 of combining the two methods. In particular, we

1 study the sub-band approach for speech recognition
2 involving partial unknown frequency-band corruption.

3

4 As a system paradigm for dealing with partial
5 frequency-band corruption, the sub-band based
6 approach has aroused much research interest over the
7 past years. In this approach, the full speech
8 frequency band is divided into several sub-bands,
9 and each sub-band is featured independently of the
10 other sub-bands, so that the local distortions in
11 the frequency band will not spread over the entire
12 feature space. Therefore, instead of requiring a
13 detailed knowledge of the noise for clearing the
14 corrupted sub-band features, the sub-band method,
15 and in general the missing feature methods, require
16 only a labelling of every sub-band/feature as
17 reliable or corrupt, for removing the unreliable
18 features from recognition.

19

20 Unfortunately, locating the corrupted sub-bands
21 itself can be a difficult task, if there is no prior
22 information on the noise. Mistakes in labelling the
23 sub-bands can cause either a loss of reliable
24 information, or an inclusion of unreliable
25 information in the recognition process. This
26 problem, i.e. extracting reliable features from a
27 sub-band observations while assuming no prior
28 knowledge on the noise, has been referred to as sub-
29 band combination.

30

31 Recent studies have suggested several methods.
32 Typically, these include the weighted-average

1 method, the neural-network method and the full-
2 combination method.

3
4 In the weighted-average method, the likelihood from
5 the individual sub-bands are combined by using a
6 geometric or arithmetic average; the contribution of
7 each sub-band is weighted by the local signal-to-
8 noise ratio (SNR) related to that sub-band.

9
10 In the neutral-net method, independent networks are
11 trained to estimate the probabilities of all
12 possible combinations of subsets of the sub-bands,
13 assuming that there exists at least one combination
14 that accounts for the clean speech. This method
15 faces the problem of how to select the best
16 combination from all the combinations given no
17 knowledge about the noisy bands. Some heuristic
18 methods, such as majority voting or distance
19 pruning, have been studied for this purpose.

20
21 The idea of explicitly creating all possible
22 combinations among the sub-bands has been further
23 studied in the full-combination model, in which the
24 likelihood of different combinations of different
25 sub-bands are combined using a weighted-average
26 method, with each weight proportional to the
27 relative reliability of a specific set of sub-bands.

28
29 In addition, the mixture of experts theory has also
30 been discussed as a possible means of sub-band
31 combination.

32

1 Clearly, a reliable estimation of the local noise
2 characteristic or SNR is crucial to the success of
3 the weighted-average model and full-combination
4 model. In fact, it is crucial to the success of all
5 missing feature methods which rely on an accurate
6 mask for labelling the reliable and corrupt regions
7 over the temporal-spectral feature space. The local
8 SNR at each time-frequency location may be estimated
9 by using the traditional spectral estimation
10 approach, involving a running estimate of the local
11 noise spectrum via spectral subtraction. This
12 method performs well when the corrupting noise is
13 stationary. But it may fail to produce accurate
14 estimates in non-stationary noise or unknown noise,
15 as in these conditions the assumption required for
16 spectral subtraction is invalidated. To overcome
17 this problem, it has been suggested that some
18 characteristics of the speech signal itself, such as
19 the harmonic nature of voiced speech may be
20 exploited for identifying the corrupted time-
21 frequency regions.

22
23 According to the present invention there is provided
24 a method of interpreting features for signal
25 processing and pattern recognition as described in
26 the attached Claims.

27
28 The present invention proposed a new approach, the
29 probabilistic union model, for combining the sub-
30 band features with unknown, time-varying partial
31 corruption. Unlike the missing feature method, the
32 new model does not require the identity of the

1 corrupted bands, instead, it combines the sub-band
2 features based on the probability theory for the
3 union of random events, to account for any possible
4 partial corruption with the sub-bands. This model
5 improves upon the previous methods in that it offers
6 robustness against partial frequency-band
7 corruption, while requiring little or no information
8 about the noise. We have incorporated the new union
9 model into an HMM framework and tested it on a
10 number of isolated word databases. The results have
11 indicated the advantage of the union model over the
12 previous methods for sub-band combination,
13 particularly for dealing with band-selective noise
14 with an unknown or time varying band location and/or
15 bandwidth.

16

17 The present invention will now be described by way
18 of example only, with reference to the accompanying
19 tables and drawings in which;

20

21 Tables I and II show experimental results showing
22 the performance of first and second embodiments of
23 the present invention against a conventional
24 technique, for incorrupt and corrupt signals
25 respectively;

26

27 Tables III and IV show experimental results showing
28 the performance of the second embodiment of the
29 present invention against further conventional
30 techniques, for stationary and non-stationary
31 corruption respectively;

32

1 Table V shows experimental results showing the
2 performance of a third embodiment of the present
3 invention in comparison to the first and second
4 embodiments and further conventional techniques;

5
6 Table VI shows experimental results showing the
7 performance of the third embodiment of the present
8 invention in comparison to the second embodiment;

9
10 Table VII shows experimental results showing the
11 performance of a fourth embodiment of the present
12 invention;

13
14 Fig. 1 illustrates the performance of a specific
15 aspect of the present invention, and;

16 Fig. 2 illustrates the raw data used to test the
17 performance of the present invention.

18

19 PROBABILISTIC UNION MODEL

20

21 A. Background

22

23 Assume a recognition system with N sub-bands, in
24 which a speech utterance may be represented by N
25 sub-band feature streams o_1, o_2, \dots, o_N , where o_n
26 represent the feature stream from the n 'th sub-band.

27 The presence of a band-selective noise can cause

28 some of the o_n 's to be corrupted. Thus, in

29 recognition we face the problem of how to extract
30 information for the utterance from a sub-band

31 feature set $\{o_1, o_2, \dots, o_N\}$, in which some of the

1 sub-band features o_n 's may be noisy, but without
2 knowledge about their identity.

3

4 When there is no noise the traditional approach for
5 extracting the information is to combine the sub-
6 band features by using the "and" (i.e. conjunction)
7 operator \wedge (although this is not usually explicitly
8 stated), i.e.

9

$$10 \quad O_{\wedge} = o_1 \wedge o_2 \wedge \dots \wedge o_N$$

11 (1)

12 where O_{\wedge} represents the combined observation.
13 Assuming that the sub-band features are independent
14 of one another, then the likelihood of O_{\wedge} , $p(O_{\wedge})$,
15 equals the product of the individual sub-band
16 likelihoods $p(o_n)$'s i.e.

$$17 \quad p(O_{\wedge}) = p(o_1 \wedge o_2 \wedge \dots \wedge o_N)$$
$$18 \quad = p(o_1)p(o_2)\dots p(o_N)$$

19 (2)

20 For convenience, we call (1) the *product model*.
21 Assume that the model, consisting of the probability
22 densities of the individual sub-bands, $P(x_n)$'s is
23 trained on clean speech to maximise the likelihood
24 of some clean utterances. When this model is used
25 for an utterance with some noisy sub-bands, then the
26 corresponding $P(o_n)$'s for the noisy o_n 's will become
27 problematic, especially if the noise is strong.
28 Typically, these noisy likelihoods may become very
29 small on the correct model because of the poor match
30 between the model and data. These small and random

1 sub-band likelihoods may easily dominate the
2 product, and then destroy the model's ability to
3 produce high likelihoods for correct phonetic
4 classes. Simply removing the sub-band likelihoods
5 with small values from the models may not improve
6 this, because low likelihoods may also be the result
7 of a phonetic mismatch, and because the likelihoods
8 corresponding to the noisy sub-bands may not be
9 small on the incorrect models which accidentally
10 match the noisy data. This problem can be improved
11 if the noisy sub-bands can be identified, whereby
12 the corresponding likelihoods can be removed or
13 "integrated" from the product, i.e. the missing
14 feature method. This identification requires the
15 local SNR related to each sub-band. This
16 information may not be available for applications
17 involving unknown, time-varying noise. This problem
18 has been addressed by using a back-off model, in
19 which each observation probability density is formed
20 as a weighted combination of two densities: one from
21 the training data and another, a uniform
22 distribution, to account for possible outliers
23 arising from the noise.

24

25 In the following we describe the probabilistic union
26 model as an alternative, to overcome the above
27 mentioned problems. We start to describe the model
28 without considering the number of noisy sub-bands
29 (except that the corruption is partial within the
30 sub-bands); then we move to an extended model which
31 takes into account knowledge on the number of noisy
32 sub-bands.

1 B. General union model

2

3 Given no knowledge about the noisy sub-bands, we can
 4 alternatively assume that, in a given set of sub-
 5 band features $\{o_1, o_2, \dots, o_N\}$, the reliable features
 6 that characterize the speech utterance may be any of
 7 the o_n 's $n = 1, \dots, N$, or any of the combinations among
 8 the o_n 's up to the complete feature set. This can be
 9 expressed, using the inclusive "or" (i.e.
 10 disjunction) operator \vee , as

11

$$12 \quad O \vee = o_1 \vee o_2 \vee \dots \vee o_N$$

13

$$= \bigvee_{n=1}^N o_n$$

14

(3)

15 where O_{\vee} is a combined observation based on \vee ,
 16 representing the reliable features within
 17 $\{o_1, o_2, \dots, o_N\}$.

18

19 For example, using a 3-band model, the expression
 20 $O \vee = o_1 \vee o_2 \vee o_3$ based on inclusive "or" assumes that
 21 the reliable features within the given $\{o_1, o_2, o_3\}$ may
 22 be o_1 , or o_2 , or o_3 , or $o_1 \wedge o_2$, or $o_1 \wedge o_3$, or $o_2 \wedge o_3$, or $o_1 \wedge o_2 \wedge o_3$.
 23 These feature combinations can characterize,
 24 respectively, a speech utterance in which there are
 25 two-band, one-band and no band corruption, therefore
 26 covering all possible partial corruptions, including

1 the no corruption case which may be encountered in a
 2 3-band system. In general, if an observation
 3 consists of N features o_1, o_2, \dots, o_N , and these features
 4 may be subjected to some partial corruption with
 5 unknown characteristics, i.e. number and location of
 6 the corrupted features and statistics of the
 7 corrupting noise, then the useful information
 8 contained in the observation may be modelled by (3).
 9 This model takes into account all possible partial
 10 corruptions, thereby requiring no knowledge on the
 11 actual corrupting noise.

12

13 If we assume that the o_n 's are discrete random
 14 vectors, then $O \vee$ is the union of the random events
 15 o_n 's. Thus, we can compute the probability $P(O \vee)$
 16 based on the rules of probability for the union of
 17 random events. This probability, for each modeled
 18 phonetic class, can then be used to decide the
 19 recognition result based on the maximum-likelihood
 20 principle. Note that $\bigvee_{n=1}^m o_n = (\bigvee_{n=1}^{m-1} o_n) \vee o_m$, so $P(O \vee)$ can
 21 be computed using a recursion

22

$$23 \quad P\left(\bigvee_{n=1}^m o_n\right) = P\left(\bigvee_{n=1}^{m-1} o_n\right) + P(o_m) - P\left(\left(\bigvee_{n=1}^{m-1} o_n\right) \wedge o_m\right)$$

24 (4)

25 for $m=2, \dots, N$. With the assumption that the o_n 's
 26 are mutually independent, then (4) can be simplified
 27 as

28

$$29 \quad P\left(\bigvee_{n=1}^m o_n\right) = P\left(\bigvee_{n=1}^{m-1} o_n\right) + P(o_m) - P\left(\bigvee_{n=1}^{m-1} o_n\right)P(o_m)$$

30 (5)

1 This computation requires only the probability
2 distributions of the individual sub-bands, i.e.
3 $P(x_n)$'s which are assumed to be estimated from *clean*
4 training data. We call (3)-(5) the *probabilistic-*
5 *union model*, which extracts information based on the
6 union of events. This is opposed to the product
7 model (1)-(2), which extracts information based on
8 the intersection of events.

9
10 Since the $P(o_n)$'s are generally not large, (5) is
11 effectively the sum of the individual sub-band
12 probabilities. A major difference between (5) and
13 (2) (i.e. the product model) is that a small $P(o_n)$
14 makes only a small contribution to (5). Therefore a
15 noisy sub-band, typically with low probability on
16 the correct model, will have little effect on the
17 union probability $P(O_v)$ associated with the correct
18 model. In other words, the union probability
19 $P(O_v)$ associated with the correct model is dominated
20 by noiseless sub-bands, unlike the product model in
21 which the likelihood associated with the correct
22 model may be dominated by those small, random and
23 noisy sub-band likelihoods. This effectively
24 increases the probability associated with the
25 correct model, such that, as long as the remaining
26 clean sub-bands contain sufficient discriminative
27 information, the correct model should still be able
28 to score highly among the competitive models.
29
30 However, (5) has a disadvantage, i.e., it
31 effectively averages the ability of each sub-band to

1 discriminate between correct and incorrect phonetic
2 classes, unlike the product model in which each sub-
3 band reinforces the other as the joint probability
4 of the sub-band features is modeled.

5

6 This characteristic makes (5) an ineffective model
7 both for utterances with more than one clean sub-
8 band, and for clean utterance without band
9 corruption. This problem may be overcome by
10 combining the use of "and" and "or" operators,
11 assuming a knowledge on the number of corrupted sub-
12 bands. This is the extended union model described
13 below.

14

15 *C. Extended Union Model*

16

17 In a first embodiment of the present invention, we
18 aim to include all the clean sub-band features into
19 a conjunction (i.e. combining them using the "and"
20 operator), such that a joint probability of the
21 clean features can be derived, which should be more
22 powerful than any of their marginal probabilities in
23 terms of discrimination. This can be achieved by
24 combining the use of "and" and "or" operators,
25 assuming only a knowledge on the number (not the
26 location) of corrupted sub-bands. Specifically, for
27 a given set of sub-band features $\{o_1, o_2, \dots, o_N\}$ if the
28 number of corrupted bands is M ($M < N$), then we know
29 that there exists one subset of $(N - M)$ sub-band
30 features which are affected little by noise. These
31 features should then be combined with the "and"
32 operator. Without knowing where the noise occurs,

1 this subset may be any of the subsets of $(N - M)$
 2 sub-band features. This uncertainty can then be
 3 modelled with the "or" operator. Combining the two
 4 together we obtain a model for representing the
 5 useful information within the given feature set

6

7

$$O \vee = \bigvee_{n_1 n_2 \dots n_{N-M}} o_{n1} o_{n2} \dots o_{n_{N-M}}$$

8

(6)

9 where the "and" operator \wedge between the o_n 's has been
 10 omitted, and the "or" operator \vee is taken over all
 11 possible combinations of N different values $(1, \dots,$
 12 $N)$ taken $(N - M)$ at a time, giving a total of ${}^N C_{N-M}$
 13 combinations. For example, in the simple case with
 14 four sub-bands, (6) can take one of the following
 15 four possible forms, corresponding to $M = 0, 1, 2$
 16 and 3, respectively:

17

18 0) $o_1 o_2 o_3 o_4$

19 1) $o_1 o_2 o_3 \vee o_1 o_2 o_4 \vee o_1 o_3 o_4 \vee o_2 o_3 o_4$

20 2) $o_1 o_2 \vee o_1 o_3 \vee o_1 o_4 \vee o_2 o_3 \vee o_2 o_4 \vee o_3 o_4$

21 3) $o_1 \vee o_2 \vee o_3 \vee o_4$

22

23 Forms 0 and 3 correspond to the product model (1)
 24 and the general union model (3), respectively, and
 25 forms 1 and 2 correspond to the assumptions that
 26 there is one and two noisy sub-bands, respectively.
 27 In form 1, for example, the union of the four
 28 conjunctions will include one conjunction providing
 29 the joint probability of all three clean sub-bands;
 30 the other three conjunctions each contain a noisy

1 sub-band, with a correspondingly low probability on
2 the correct model, and therefore make only a small
3 contribution to the union probability associated
4 with the correct model. In a similar way, in form 2
5 assuming two noisy sub-bands, one of the six
6 conjunctions will correspond to the remaining two
7 clean sub-bands and this conjunction will dominate
8 the union probability associated with the correct
9 model.

10

11 For convenience, we call (6) a *union model of order*
12 *M*. As indicated above, the value of *M* corresponds
13 to the maximum number of noisy sub-bands that can be
14 accommodated in the model, in terms of leaving at
15 least one conjunction consisting of only clean sub-
16 bands. The product model (1), which includes a full
17 conjunction of the sub-bands, corresponds to a union
18 model with order $M = 0$ and therefore is best
19 suitable for clean utterance without band
20 corruption. The general union model (3) has an
21 order $M = N - 1$, and thus may accommodate up to $N -$
22 1 noisy bands. Note that while a match between the
23 order of the model and the number of noisy bands is
24 desirable to maximise the information being
25 extracted, a union model with order *M* may also be
26 suited to situations where the number of noisy sub-
27 bands is less than *M*.

28

29 For example, the above form 2, with order $M = 2$, may
30 also be used to accommodate one noisy sub-band or
31 none. This offers robustness against uncertainty on
32 the number of corrupted bands. This characteristic

1 has been exploited previously for the selection of
2 the model order, to seek a balance between the
3 maximum performance and robustness. Details of this
4 will be discussed later, along with a new algorithm
5 for automatic order selection.

6

7 The expression for the union probability of (6) can
8 be readily derived with o_n in (5) replaced by the
9 appropriate conjunctions of sub-band features, i.e.

10 $o_{n1} o_{n2} \dots o_{nN-M}$, assuming independence between the
11 features. This computation requires only the
12 probability distributions of the individual sub-
13 bands, as was required in the general union model
14 discussed in the previous section.

15

16

IMPLEMENTATION

17

18 In this section, we first describe the
19 implementation of the union model within a HMM
20 framework, and then we describe the algorithms
21 proposed for order selection.

22

23 A. Incorporation into HMM

24

25 We have built the above union model (6) into an HMM
26 for combining the sub-band features at the *frame*
27 level. Assume that there are N sub-bands, and that
28 a speech utterance in each sub-band is represented
29 by a sequence of frame vectors $o_n(1), o_n(2), \dots$

30 $o_n(T), n = 1, \dots, N.$

31

1 Combining the sub-band features at the frame level
 2 means that the union model (6) is applied at every
 3 frame time t , to combine the frame vectors $o_1(t)$,
 4 $o_2(t)$, ... $o_N(t)$ from all the sub-bands to obtain a
 5 union observation $O_v(t)$, $t = 1, \dots, T$. Then we
 6 modify the conventional HMM for this new observation
 7 sequence, by using a union-based observation
 8 probability distribution for each $O_v(t)$. This HMM
 9 can be written as

10

11

$$P(O|\lambda) = \sum_S P(S|\lambda) \prod_{t=1}^T B_{s_t}(O_v(t))$$

12

(7)

13 where O represents the frame sequence for all the
 14 sub-bands, $P(S|\lambda)$ is the probability of the state
 15 sequence S , and $B_i(O_v)$ is the union based frame-
 16 level observation probability distribution in state
 17 i . As usual, the parameter set of the model, λ ,
 18 includes the state transition probability matrix and
 19 initial state probability vector, which are needed
 20 for calculating the probability $P(S|\lambda)$ and the
 21 observation distribution set $\{B_i(O_v)\}$. As described
 22 above with the assumption that the sub-band frames
 23 are mutually independent, the probability $B_i(O_v)$ is
 24 only a function of the individual probabilities
 25 $B_i(o_n)$'s where $B_i(o_n)$ represents the observation
 26 probability of the frame in sub-band n and state i .
 27 For a discrete-observation HMM, these sub-band
 28 observation probability distributions are readily
 29 available, and so $B_i(O_v)$ can be readily calculated

1 by using the algorithm described above. However,
 2 note that (4) or (5), for computing the union
 3 probability, apply only to probabilities, not to
 4 probability densities or likelihoods. Therefore a
 5 special treatment is needed to resolve this issue
 6 when implementing the union model for a continuous-
 7 observation HMM, which employs an observation
 8 probability density $b_i(o_n)$ to account for the frame
 9 in sub-band n and state i . Basically, we seek an
 10 approximated probability based on a likelihood.
 11 However, this approximation is not needed in the
 12 model training stage, if the model is trained on
 13 clean speech data. Although $B_i(Ov)$ varies with the
 14 order M for recognition, there is only one form,
 15 with order $M = 0$, that best matches a clean
 16 observation. Therefore in the training stage we can
 17 compute the union observation probability $B_i(Ov)$ as
 18 the full conjunction probability $B_i(o_1)...B_i(o_N)^1$. Since
 19 this probability is proportional to the likelihood
 20 $b_i(o_1)...b_i(o_N)$, we can train the model by maximising
 21 the likelihood function

22

$$23 \quad p(O|\lambda) = \sum_S P(S|\lambda) \prod_{t=1}^T \prod_{n=1}^N b_{s_t}(o_n(t))$$

24 (8)

25 ¹ More rigorously, the probability of a
 26 continuous o_n should be written as $B_i(x \in \Omega_n)$
 27 i.e. the probability of a continuous random
 28 vector x falling into a sub-space Ω_n

1 surrounding o_n . But for simplicity we will
 2 keep using the expression $B_i(o_n)$.
 3
 4 and this can be accomplished by using the
 5 standard forward-backward re-estimation
 6 algorithm. In recognition, decisions are made
 7 by comparing the probability $P(O|\lambda)$, defined in
 8 (7), between different models. As with the
 9 conventional HMM, this probability can be
 10 computed by using the Viterbi algorithm, i.e.

$$12 \quad \delta_t(j) = \max_i (\delta_{t-1}(i) + \log a_{ij}) + \log B_j(O_v(t))$$

13 (9)

14 where $\delta_t(i)$ is the log probability associated
 15 with a best state-sequence ending in state i
 16 for the observation up to time t , and a_{ij} is the
 17 state transition probability. With order M
 18 $\neq 0$, there may be two ways to obtain an
 19 approximated union probability $B_i(O_v)$, based on
 20 the sub-band frame likelihoods $b_i(o_1), \dots, b_i(o_N)$.
 21 One way is to leave out the product term in
 22 (5), assuming that it is small and can be
 23 neglected in comparison to the other two
 24 additive terms. As such, the union probability
 25 $B_i(O_v)$ with O_v defined by (6) can be written as
 26

$$27 \quad B_i(O_v) \cong \sum_{n_1 n_2 \dots n_{N-M}} B(o_{n_1}) B_i(o_{n_2}) \dots B_i(o_{n_{N-M}})$$

$$\propto \sum_{n_1 n_2 \dots n_{N-M}} b_i(o_{n_1}) b_i(o_{n_2}) \dots b_i(o_{n_{N-M}}) \quad (10)$$

1 where the summation is over all possible
2 combinations of N different values $(1, \dots, N)$ taken
3 $(N - M)$ at a time. Therefore (10) indicates a
4 likelihood that may be used to approximate the union
5 probability.

6
7 Alternatively, a sigmoid function may be used to
8 approximate the sub-band frame probability $B_i(o_n)$
9 based on the likelihood $b_i(o_n)$, i.e.

10

$$11 \quad B_i(o_n) \cong \frac{1}{1 + e^{-\ln b_i(o_n)}} \quad (11)$$

12

13 This has the property that it produces an
14 approximated probability that is proportional to the
15 likelihood value, and at the same time satisfies the
16 constraint $0 \leq B(o_n) < 1$ (this is required by (5) not
17 to produce a negative probability). The probability
18 $B_i(O_v)$ with each $B_i(o_n)$ defined by (11) can thus be
19 computed based on (5), including the product term.
20 Because this term is usually very small
21 (particularly for models with an order $M \ll N$), the
22 two methods described above are based on (10) and
23 (11) have been found to produce almost identical
24 results.

25

26 Based on the assumption that the conjunction
27 including only the clean bands should dominate the
28 union probability for the correct model, (10) may be
29 further approximated as

30

$$B_i(O_v) \cong \max_{n_1 n_2 \dots n_{N-M}} b_i(o_{n_1}) b_i(o_{n_2}) \dots b_i(o_{n_{N-M}}) \quad (12)$$

where the maximisation is over all possible combinations of N different values $(1, \dots, N)$ taken $(N - M)$ at a time. We have found in our experiments that, given the same order M ($M > 0$), the recognition results base on (10) and (12) are similar for low SNR conditions. However, in high SNR conditions, (10) was usually found to perform significantly better than (12). This is because (10) does not physically remove any sub-bands from recognition which (12) does. In high SNR conditions, those bands thrown away in (12) may still carry useful information.

B. Algorithms for order selection

A second embodiment of the present invention enables selection of an appropriate order to accommodate the corrupted sub-bands within an observation. As indicated above if there is no knowledge on the corrupting noise, it is safer to select a high order to accommodate as much noise as possible. However, because a higher-than-needed order will usually cause a loss of information due to unnecessary disjunction of the clean sub-bands, the order must be subject, for example, to an acceptable performance for clean speech recognition. We call this the balance fixed-order algorithm, which has been tested previously and has shown a limited success. In the following we describe an improved

1 algorithm, which derives the order automatically
2 based on an optimality criterion.

3

4 As discussed above, an overestimated order (i.e. an
5 order larger than the actual number of corrupted
6 sub-bands) will lead to an unnecessary disjunction
7 between the clean bands. This can cause some of the
8 information relating to the joint probability
9 distribution of the clean bands to be lost. On the
10 other hand, an underestimated order (i.e. an order
11 smaller than the actual number of corrupted sub-
12 bands) will cause every conjunction in the union
13 model to include, and so to be affected by, one or
14 more corrupted sub-bands. Formally, we define the
15 matched order as the order that equals the number of
16 corrupted sub-bands. With this order, the union
17 model will include a conjunction which contains all
18 of the clean sub-bands together and no others,
19 thereby capturing more discriminative information
20 than either of the order-overestimated model or
21 order-underestimated model, i.e. the order
22 mismatched model. Because the order-matched model
23 captures more clean band information, it should have
24 more characteristics of a clean utterance than the
25 order-mismatched model. This assumption forms the
26 basis of our order selection algorithm. In
27 particular, we use the state duration probability
28 for clean utterance to estimate the matched order.

29

30 The state duration probability $P_i^u(d)$, for d frames
31 in state i of phonetic unit u , is estimated in the
32 training stage using the clean training data. Given

1 training stage using the clean training data. Given
 2 a test utterance, we perform recognition by using a
 3 set of union models, each with a different order,
 4 assuming that these will include the matched order.
 5 For each order, we obtain a recognition result (in
 6 the form of a unit sequence) $U(r) = u_1(r)u_2(r)\dots u_n(r)$
 7 where r is the order index, along with the
 8 associated state duration $d_i(r)$, for each state i of
 9 $U(r)$. Because the model with the matched order
 10 captures the maximum clean band information, its
 11 state duration should be most similar to the state
 12 duration of a clean utterance. Therefore an
 13 appropriate estimate of the matched order would be
 14 the order whose associated state duration has the
 15 maximum probability, i.e.

16

$$\hat{r} = \arg \max_r \frac{1}{S(r)} \sum_{u \in U(r)} \sum_{i \in u} \ln P_i^u(d_i(r))$$

(13)

18

19 where $S(r)$ stands for the total number of states in
 20 $U(r)$. The final recognition result is then given by
 21 $U(\hat{r})$.

22

23

EXPERIMENTS

24

25 The TIDIGITS connected digit database was used to
 26 evaluate the performance of the new union model.
 27 This database contained connected digit strings from
 28 225 adult speakers, conveniently divided into
 29 training and testing sets. The testing set
 30 contained a total of 6196 utterances from 113

1 speakers, each speaker contributing five utterances,
2 containing 2, 3, 4, 5 and 7 digits, respectively.
3 In recognition we assumed no advance knowledge of
4 the number of digits in an utterance.

5
6 The speech was sampled at 8 kHz, and was divided
7 into frames of 256 samples, with a between-frame
8 overlap of 128 samples. For each frame, we used a
9 mel-scaled filter bank to estimate the log-amplitude
10 spectra of speech. Based on these log filter-bank
11 spectra, both the full-band features and sub-bands
12 features were calculated. The full-band features
13 were used for comparison, which were the full-band
14 MFCCs (mel-frequency cepstral coefficients) and were
15 obtained by taking a DCT over the complete set of
16 the log filter-bank spectra. The sub-band features
17 were obtained by first grouping the filter-bank
18 channels uniformly into sub-bands, and then, for
19 each sub-band, performing a DCT for the log filter-
20 bank spectra within that sub-band. This gives the
21 sub-band MFCCs. In both cases, the first-order
22 delta MFCCs were included in the feature vectors.
23 The division of the speech frequency-band into sub-
24 bands remains a topic of research. To effectively
25 isolate any local frequency corruption from the
26 other usable bands, a fine subdivision may be
27 desirable. However, breaking the available
28 frequency-band into too many independent sub-bands
29 will cause much of the spectral dependency to be
30 ignored, thus giving a poor phonetic discrimination.
31 As an experimental study, we have tested the
32 division of the available frequency-band into 3, 5

1 and 7 sub-bands, respectively, earlier for the E-set
2 word recognition and now for the connected digit
3 recognition. Both experiments indicate that the 5-
4 band model appears to be a better choice in terms of
5 the balance between the noise localisation and
6 phonetic discrimination. Therefore in the following
7 we focus on the experiments with five sub-bands
8 (i.e. $N = 5$, in models (2) and (6)).

9
10 Specifically, these five sub-bands were grouped from
11 a mel-scaled filter bank with 30 channels, each sub-
12 band thus containing six log filter-bank spectral
13 components for a frame. From these six components
14 three MFCCs were derived, plus the delta parameters,
15 as the feature vector of a sub-band frame. Thus,
16 for this 5-band system, the overall size of the
17 feature vector for a frame is $5 \times 6 = 30$. The full-
18 band feature vector of a frame includes 20
19 components (10 MFCCs and 10 delta MFCCs), derived
20 from a mel-scaled filter bank with 20 channels.

21
22 In addition to the union model, for comparison, we
23 also implemented a baseline HMM which used the above
24 full-band features and a product model which is a
25 special case of the union model with order $M = 0$.
26 All these models were based on Gaussian mixture
27 densities with diagonal covariance matrices, and
28 were trained on clean training data. In particular,
29 each digit was modelled with 10 states, and a
30 silence model with one state was built to account
31 for the silences surrounding each utterance and the
32 optional silences between digits. Each of these

1 states contained eight mixtures. For the union
2 model, we also recorded the histograms of state
3 occupancy occurring in each digit, as the estimates
4 of the state duration probabilities. The state
5 duration probability was used only for selecting the
6 model order, as described above and was not
7 incorporated into the HMMs for scoring the
8 observations.

9
10 In the following we first present the recognition
11 results by the union model under various testing
12 conditions. Then we discuss its generalisation to
13 the combination of different types of feature
14 streams, and its combination with a conventional
15 noise-reduction technique.

16 17 *A. Tests with clean speech*

18
19 Table I presents the string accuracy obtained by the
20 union model and the baseline model, respectively,
21 for clean utterance recognition. As shown in the
22 table, our baseline HMM achieved a string accuracy
23 of 97.53%,

24 Based on (6), for a union model with N sub-bands
25 (now $N = 5$), recognition can be performed with
26 different orders (i.e. M) within the range
27 $0 \leq M \leq N-1$ (now $0 \leq M \leq 4$). Table I presents the
28 accuracy obtained by using each of these individual
29 orders, along with the accuracy based on the
30 automatically selected order. Note that at order
31 0, the union model is equivalent to a product
32 model.

1 As described earlier, since there is no band
2 corruption, a clean speech utterance is better
3 characterised by a full conjunction of all the sub-
4 band features. This explains why the product model,
5 derived from such a conjunction, produced the best
6 performance among all the orders within the range
7 $0 \leq M \leq 4$. As expected, the performance of the union
8 model decreased as the order was increased, because
9 of the disjunction between the clean sub-band
10 features.

11
12 Given a test utterance, the above models with fixed
13 orders each produced a recognition result, tagged by
14 the associated order. The automatic order selection
15 algorithm, (12), was then applied to these results
16 to select an order with maximum state duration
17 probability, thereby obtaining the final recognition
18 result. As shown in Table I, this gives an accuracy
19 that is very close to the accuracy obtained by the
20 best (i.e. matched) order - order 0. Fig. 1 shows
21 the histograms of the orders selected by the
22 algorithm. As indicated in Fig. 1(a), for clean
23 test utterances, the algorithm correctly selected
24 more than 50% of the orders. This correct selection
25 rate may be improved by putting a restriction on the
26 order range searched by the algorithm. For example,
27 we tested the use of a smaller range $0 \leq M \leq 3$
28 instead of $0 \leq M \leq 4$ and ended with slightly better
29 result for clean utterance recognition. However,
30 allowing the uncertainty of the environment, in the
31 following all automatic orders were selected from
32 the order range $0 \leq M \leq 4$.

1 *B. Tests with stationary band-selective noise*

2

3 To evaluate the robustness of the union model, we
4 first tested the model for the utterances corrupted
5 by stationary band-selective noise. The noise,
6 added to the speech, was generated by passing
7 Gaussian white noise through a band-pass filter with
8 a 3-dB cut-off bandwidth of 100 Hz and a varying
9 central frequency. In particular, six different
10 central frequencies were considered, these were 600
11 Hz, 850 Hz, 1200 Hz, 1500 Hz, 2000 Hz and 2500 Hz.
12 These were chosen to create the effects that there
13 were one sub-band, two sub-band and three sub-band
14 corruptions, respectively, within the five sub-bands
15 of the system. Specifically, the noises with
16 central frequencies 600 Hz, 1200 Hz and 2000 Hz were
17 located within sub-band 2, 3 and 4, respectively,
18 and each thus caused only one sub-band corruption;
19 the noises with central frequencies 850 Hz, 1500 Hz
20 and 2500 Hz were located around the border of sub-
21 bands 2 and 3, 3 and 4, and 4 and 5, respectively,
22 and each thus caused two sub-band corruptions. The
23 noises corrupting three sub-bands were generated by
24 combining two noise components with different
25 central frequencies, in particular, 600 Hz and 1500
26 Hz (corrupting sub-bands 2, 3 and 4), and 1200 Hz
27 and 2500 Hz (corrupting sub-bands 3, 4 and 5),
28 respectively. The six band-selective noises, plus
29 the two combined noises, resulted in a total of
30 eight different noise conditions. For all
31 conditions, we assumed no prior knowledge of the
32 noise being available for the union model.

1 Table II presents the recognition results, as a
2 function of the number of corrupted sub-bands and
3 SNR within each test utterance. These results are
4 averaged over the appropriate noise conditions
5 producing the same number of noisy sub-bands, as
6 elaborated above. From Table II, two particularly
7 useful observations can be made for the union model.
8 Firstly, for each given SNR condition, the fixed-
9 order model achieved the maximum accuracy at the
10 order that matched the number of corrupted sub-
11 bands. Secondly, the automatic-order model was able
12 to achieve an accuracy that was close to the
13 matched-order accuracy, throughout all test
14 conditions. In particular, we see that in two cases
15 (with three noisy bands, SNR=10 dB and 5 dB,
16 respectively) the automatic-order model achieved a
17 higher recognition accuracy than the corresponding
18 matched-order accuracy (i.e., 76.27% vs 72.32%, and
19 64.81% vs 64.75%, respectively). This may be
20 because the order selection algorithm is operated on
21 each utterance basis, so it may choose an order
22 which includes some noisy bands, in which the local
23 SNRs are high. Fig. 1(b)-(d) show the histograms of
24 the orders selected by the algorithm for the noisy
25 conditions. We see that in each condition, the
26 algorithm selected the matched order with the
27 highest frequency. Based on Tables I and II, we
28 then may conclude that, equipped with the automatic
29 order selection algorithm, the union model can
30 effectively achieve a near matched-order performance
31 for both clean and noisy conditions, without
32 requiring any information on the nature of the

1 environment (i.e. clean or noisy) and on the noise
2 (i.e. the location and number of noisy sub-bands),
3 if the environment is noisy.

4
5 We next conducted comparisons between the union
6 model with automatic order and hence requiring no
7 knowledge on the noise, with two other models with
8 knowledge on the noise. The first model we compared
9 was an ideal missing-feature model, or the "oracle"
10 model which assumed a full a priori knowledge of the
11 corrupted sub-bands and removed those bands manually
12 from the recognition. The second model being
13 compared was a baseline HMM equipped with a Wiener
14 filtering front-end for removing the noise, based on
15 the assumption that the noise was stationary and for
16 which a spectral estimate was available. The
17 spectrum of the stationary band-selective noise was
18 estimated in the interval without speech. The
19 spectral estimate was then used to build a Wiener
20 filter, derived from spectral subtraction to enhance
21 the noisy signal before recognition. Table III
22 presents the results. As expected, the oracle model
23 performed better than the union model, and the gap
24 between their performances is significant in many
25 cases. Later we will discuss an improvement over
26 the union model, to reduce this performance gap. In
27 one case, with three noisy sub-bands and SNR=10 dB,
28 the union model outperformed the oracle model. This
29 is because throwing away the three bands with
30 relatively high SNR in the oracle model caused a
31 loss of much useful information. However, when all
32 these bands were included, it gave an accuracy of

1 only 28.18%, as shown in Table I. So a "soft"
2 rather than a binary decision is preferred as to
3 whether to include or exclude a particular sub-band.
4 The union model provides such a soft-decision
5 mechanism. It is capable of ignoring those noisy
6 bands that significantly violate the statistics of
7 the training data population; but it physically
8 removes no band from recognition, as such each band
9 retains a contribution, proportional to its
10 likelihood value, to recognition. Comparing Table
11 III with Table II, we see that the Wiener filtering
12 considerably improved the performance of the
13 baseline model. However, the union model still
14 performed significantly better than the baseline
15 model with Wiener filtering, throughout all test
16 conditions.

17

18 *C. Test with real-world, non-stationary noise*

19

20 Next, we tested the union model, with automatic
21 order, for recognising utterances corrupted by some
22 real-world noises. The noise data used in the
23 experiments are shown in Fig. 2, which include the
24 sounds of a ding, a telephone ring, a whistle, which
25 were extracted from the sound files "ding.wav",
26 "ring.wav" and "whistle.wav", respectively, from the
27 Windows operating system, and the sounds of
28 "contact" and "connect", which were used in an
29 internet tool (ICQ) for on-line contact, chat and
30 sending messages. These noises each have a dominant
31 band-selective characteristic, and the noises
32 "contact" and "connect" are particularly non-

1 stationery. These noises were added, respectively,
2 to each of the test utterances for recognition
3 experiments. Table IV presents the string accuracy
4 obtained for each of these noises and the average
5 accuracy over all these noises. As a reference,
6 Table IV also includes the results given by the
7 baseline model. No noise reduction technique was
8 employed in the baseline model, due to the non-
9 stationary nature of the noise and due to the
10 assumption that there was no prior knowledge about
11 the noise.

12

13 Table IV indicates that the performance of the union
14 model for the telephone-ring noise and "connect"
15 noise is less significant in comparison to the
16 performance for the other three types of noise.
17 This is because both the telephone-ring noise and
18 "connect" noise have particular multi-band
19 characteristics. For the telephone-ring noise, for
20 example, the first two tones lay in bands 3 and 4,
21 respectively, and the last two tones fell within
22 band 5, which thus affected 3 sub-bands. We have
23 experienced weakness of the sub-band method for
24 dealing with wide-band noise. Wide-band noise
25 affects all sub-bands, which therefore violates the
26 noise-localization assumption made in the sub-band
27 model. For a system to be capable of dealing with
28 both narrow-band and wide-band noises, a combination
29 of different techniques may be needed. We will show
30 such an example later.

31

1 D. Generalisation to partial feature stream
2 corruption

3
4 So far we have described a union model for
5 extracting useful features from a set of sub-band
6 feature streams $\{o_1, o_2, \dots, o_N\}$, where each o_n
7 represents the feature stream of a specific sub-
8 band. In a third embodiment of the present
9 invention, this model may be generalised by
10 considering the feature set $\{o_1, o_2, \dots, o_N\}$, as a
11 collection of more types of feature stream rather
12 than only the sub-band feature stream. In speech
13 recognition, a speech utterance may be represented
14 by multiple feature streams, typically, the static
15 spectra and dynamic spectra, over varying time
16 scales. In real-world applications, due to the
17 background noise or channel effects, there may be
18 only a subset of the given feature streams that
19 remain reliable. For example, the static spectral
20 features are more sensitive to a stationary or
21 slowly-varying noise than the dynamic spectral
22 features. If a feature stream is adversely
23 affected, it should play a less significant role
24 than the other unaffected streams in recognition.
25 However, without prior knowledge of the
26 environmental or noise condition, it can be
27 difficult to decide which subset of the feature
28 streams provides reliable information. This
29 uncertainty may be dealt with by using the union
30 model. For this, we rephrase the above sub-band
31 combination problem as a general feature selection
32 problem, i.e. selecting reliable features from a

1 feature set $\{o_1, o_2, \dots, o_N\}$, where each o_n represents a
2 specific feature stream, given that some of the o_n 's
3 may be corrupted, but without knowledge about their
4 identity.

5
6 As an application, we have generalised our previous
7 sub-band union model by applying the union not only
8 to the combination of the sub-bands, but also to the
9 combination of the static and dynamic feature
10 streams, to further select the feature stream within
11 each sub-band that is least affected by noise.
12 Specifically, we separated the static feature and
13 dynamic feature within each sub-band into two
14 feature streams o_n and Δo_n , where Δo_n represents the
15 dynamic feature stream (i.e. Δ MFCCs), and then we
16 modelled the entire feature set $\{o_1, \dots, o_N, \Delta o_1, \dots, \Delta o_N\}$
17 with a union model with $2N$ input streams and a full
18 order range $0 \leq M \leq 2N-1$. With the previously defined
19 5-band system, we then had a union model with 10
20 input feature streams (five for MFCCs and five for
21 Δ MFCCs, each consisting of 3 components for each
22 frame) and a full range order $0 \leq M \leq 9$. Using this
23 generalised union model, we repeated all the
24 previous experiments under exactly the same test
25 conditions. The generalised model used automatic
26 orders selected from an order range $0 \leq M \leq 8$.

27
28 Tables V and VI present the string accuracy obtained
29 by the generalised union model, along with the
30 average error reduction in comparison to the
31 previous union model without applying the union for

1 the static and dynamic feature streams, as shown in
2 Tables I, III and IV. Comparing Table V with Table
3 I, we see that the generalised union model even
4 improved the accuracy for clean utterance
5 recognition. Comparing Table V with Table III, for
6 stationary band-selective noise, we see that the
7 generalised model significantly improved over the
8 previous union model for all noise conditions,
9 particularly for the conditions with multiple noisy
10 bands. Comparing Table V with the oracle model in
11 Table III, we see that the generalised union model
12 outperformed the oracle model in many cases, and it
13 actually achieved better average performance than
14 the oracle model. Table VI shows the string
15 accuracy by the generalised union model in real-
16 world, non-stationary noise, corresponding to Table
17 IV. Comparing these two tables, we again see that
18 the generalised union model significantly improved
19 the accuracy for all noise conditions. Improvements
20 for the noisy cases may be due to the separation and
21 removal of those static features that were more
22 adversely affected by the noise.

23 *E. Combination of Techniques*

24
25 So far we have assumed no prior knowledge about the
26 noise. This is typical for some random, abrupt
27 noises. However, real-world noise may be a mixture
28 of stationary noise and abrupt noise. For
29 stationary noise, with reasonably sufficient
30 observations, it is possible to obtain an estimate
31 of the noise characteristics. In a fourth
32 embodiment of the present invention, we consider the

1 building of a system in which the union model and
2 some conventional noise-reduction techniques are
3 combined, to deal with this type of mixed noise.
4 The stationary noise component may be removed, for
5 example, by spectral subtraction for additive noise,
6 or by cepstral mean subtraction for convolutive
7 noise. The remaining unknown unexpected noise
8 component can be dealt with by the union model if it
9 has a band-selective characteristic.
10
11 We have tested such a system by creating noisy
12 speech data involving both stationery noise and
13 unknown, band-selective noise, both being additive.
14 Specifically, the stationary noise was a car noise,
15 obtained from the Aurora 2 database, which exhibited
16 a wide-band characteristic; the band-selective noise
17 was a whistle, as shown in Fig. 2, which simulated a
18 further unknown and unexpected band-selective
19 corruption occurring to the utterance. To reduce
20 the stationary noise component, we may use the
21 Wiener filtering technique as described above. Here
22 we considered a different technique, i.e. noise
23 compensation. In particular, we assumed that we had
24 the models trained in the car environment, so that
25 the mismatch between the model and data, due to the
26 existence of the stationary noise, could be reduced.
27 While we assumed knowledge about the occurrence of
28 the stationary noise, we assumed no knowledge about
29 the occurrence of the whistle during the utterance.
30 The SNR's of the two noise components were
31 calculated separately relative to the clean speech
32 data, and each was 10 dB (so the overall SNR within

1 each utterance was about 7 dB). The generalised
2 union model described above was used in this
3 experiment. Table VII presents the recognition
4 results, showing the advantage of the combination of
5 the union model and noise compensation technique for
6 dealing with the mixed noise.

7
8 We then further developed this combination into a
9 simple parallel-environment model, in which two sets
10 of generalised union models, trained for clean
11 condition and car condition respectively, were run
12 in parallel, and the final result was selected using
13 the order selection algorithm over the two sets of
14 models. This model removes the requirement for a
15 knowledge of the environment (i.e. clean or car).
16 For clean speech input, this model produced a string
17 accuracy of 95.30%, and for the noisy speech input,
18 assuming the same mixed noise as described above,
19 this model produced a string accuracy of 74.66%.
20 Both accuracies were close to their respective
21 environment-model matched accuracy, i.e. 96.21% and
22 75.21%, shown in Table V and Table VII,
23 respectively.

24
25 It will be appreciated that various improvements and
26 modifications can be made without departing from the
27 scope of the invention.

28
29 Whilst the invention has been described with
30 specific embodiments relating to speech recognition,
31 it will be appreciated that the invention is
32 applicable to any other areas of signal processing

1 and pattern recognition involving partial unknown
2 feature corruption, for example, image processing,
3 statistical language processing, communication, and
4 artificial intelligence.

5

6 Alternative techniques for dealing with known or
7 trainable noise or environmental effects may be
8 incorporated into the invention, for example,
9 speaker adaptation to accommodate speaker variation,
10 or recognition of key words.

11

12 In the context of speech recognition, the principle
13 of the invention can be extended to the combination
14 of units at a higher level, for example phoneme or
15 syllable.

TABLE I

STRING ACCURACY (%) FOR CLEAN UTTERANCES, FOR THE UNION MODEL WITH FIXED ORDERS AND AUTOMATICALLY SELECTED ORDER (AO), AND FOR THE BASELINE HMM. AT ORDER 0, THE UNION MODEL IS EQUIVALENT TO A PRODUCT MODEL

Union Model						Baseline HMM
Order						
0	1	2	3	4	AO	
96.48	95.08	92.03	86.99	64.11	95.58	97.53

TABLE II

STRING ACCURACY (%) IN STATIONARY BAND-SELECTIVE NOISE, FOR THE UNION MODEL WITH FIXED ORDERS AND AUTOMATICALLY SELECTED ORDER (AO), AND FOR THE BASELINE HMM. THE MATCHED-ORDER ACCURACY FOR THE UNION MODEL IS SHOWN IN ITALIC

SNR (dB)	#	Union Model						Baseline HMM
	Corrupted Bands	Order						
		0	1	2	3	4	AO	
10	1	58.04	92.81	89.92	81.52	52.93	90.67	61.62
	2	47.33	76.47	88.65	79.11	46.85	86.63	63.16
	3	28.18	59.74	72.13	72.32	42.88	76.27	34.20
5	1	40.98	90.60	87.24	76.77	46.64	88.29	37.04
	2	31.04	61.10	86.82	76.10	42.87	83.91	38.85
	3	9.35	35.55	53.50	64.75	37.10	64.81	13.66
0	1	24.35	85.33	82.33	69.38	37.58	83.93	17.70
	2	20.05	42.94	83.95	71.35	38.20	79.89	20.32
	3	2.86	20.27	34.47	56.57	31.58	53.95	3.77

TABLE III

COMPARISONS OF STRING ACCURACY (%) IN STATIONARY
BAND-SELECTIVE NOISE, FOR THE UNION MODEL WITH
AUTOMATIC ORDER, FOR THE ORACLE MODEL WITH A FULL A
PRIORI KNOWLEDGE OF THE NOISY BANDS, AND FOR THE
BASELINE HMM WITH WIENER FILTERING (WF)

SNR (dB)	Model	#Corrupted Bands			Average
		1	2	3	
10	Union	90.67	86.63	76.27	84.52
	Oracle	94.31	89.65	66.73	83.56
	Baseline	79.61	81.81	64.42	75.28
	(WF)				
5	Union	88.29	83.91	64.81	79.00
	Oracle	93.21	88.39	65.18	82.26
	Baseline	60.47	62.15	36.95	53.19
	(WF)				
0	Union	83.93	79.89	53.95	72.59
	Oracle	89.83	86.46	62.85	79.71
	Baseline	29.13	36.76	14.61	26.83
	(WF)				

TABLE IV

STRING ACCURACY (%) IN REAL-WORLD NON-STATIONARY
NOISE, FOR THE UNION MODEL WITH AUTOMATIC ORDER, AND
FOR THE BASELINE HMM

SNR (dB)	Model	Noise Type					Average
		Ding	Tel Ring	Whistle	Contact	Connect	
10	Union	85.30	72.85	88.62	87.41	74.13	81.66
	Baseline	65.28	60.23	50.44	53.62	41.59	54.23
5	Union	80.46	60.77	86.06	84.60	58.76	74.13
	Baseline	43.56	34.49	25.87	30.57	16.03	30.10
0	Union	75.02	50.81	82.18	79.62	36.73	64.87
	Baseline	22.26	17.75	8.28	14.27	4.50	13.41

TABLE V

STRING ACCURACY (%) FOR CLEAN SPEECH AND IN
STATIONARY BAND-SELECTIVE NOISE, FOR THE GENERALISED
UNION MODEL, AND AVERAGE ERROR REDUCTION (%) IN
COMPARISON TO THE PREVIOUS UNION MODEL IN TABLES I
AND III, ALL WITH AUTOMATIC ORDERS

SNR (dB)	# Corrupted Bands			Average	Ave. Error Reduction
	1	2	3		
Clean	96.21				14.25
10	92.49	90.75	86.45	89.90	34.75
5	90.55	88.06	80.38	86.33	34.90
0	87.10	84.65	70.97	80.91	30.35

TABLE VI

STRING ACCURACY (%) IN REAL-WORLD NON-STATIONARY NOISE, FOR THE GENERALISED UNION MODEL, AND AVERAGE ERROR REDUCTION (%) IN COMPARISON TO THE PREVIOUS UNION MODEL IN TABLE IV, BOTH WITH AUTOMATIC ORDERS

SNR (dB)	Noise Type					Average	Ave. Error Reduction
	Ding	Tel Ring	Whistle	Contact	Connect		
10	90.96	81.99	90.95	88.15	79.21	86.25	25.02
5	88.12	73.87	88.75	85.31	65.80	80.37	24.12
0	84.68	62.90	84.88	81.81	44.96	71.85	19.86

TABLE VII

STRING ACCURACY (%) IN MIXED STATIONARY WIDE-BAND NOISE (CAR) AND UNKNOWN BAND-SELECTIVE NOISE (WHISTLE), EACH WITH AN SNR=10 DB, SHOWING THE EFFECTIVENESS OF COMBINING THE NOISE COMPENSATION TECHNIQUE AND THE UNION MODEL

	No Noise Compensation	With Noise Compensation
Union	35.75	75.21
Baseline	35.93	56.55

CLAIMS

1. A method of interpreting features for signal processing and pattern recognition in which recognition of a signal or pattern is enabled by a model in which the sample to be interpreted is considered as a set of N observations, M of which are corrupt, and a disjunction is performed over all possible combinations of N different values $(1, \dots, N)$ taken $N-M$ at a time.
2. A method as claimed in Claim 1 wherein $0 < M \leq N-1$.
3. A method as claimed in either preceding Claim in which the value of M , namely the number of corrupt observations, defines an order of the model, and is estimated using an optimality criterion in which:

it is assumed that the matched order is the order having the most characteristics of a clean signal,
an aspect of the clean signal is selected,
the values of the aspect are compared for different orders,
and

the chosen order is defined as the order for which the value of the aspect is closest to that of a clean signal.
4. A method as claimed in any preceding Claim wherein the signal to be processed is a speech signal.
5. A method as claimed in any preceding Claim wherein the set of N observations comprises a set of N sub-band feature streams.
6. A method as claimed in Claim 3 in which said selected aspect is a state duration probability.

7. A method as claimed in Claim 6 in which the optimality criterion is obtained from the order selection algorithm

$$\hat{r} = \arg \max_r \frac{1}{S(r)} \sum_{u \in U(r)} \sum_{i \in u} \ln P_i^u(d_i(r))$$

where: r is the order index;

\hat{r} is the order index with the highest associated state duration probability;

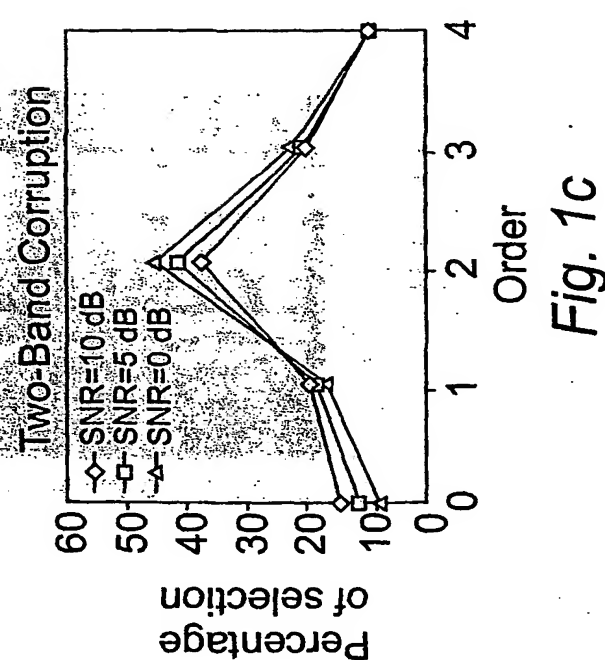
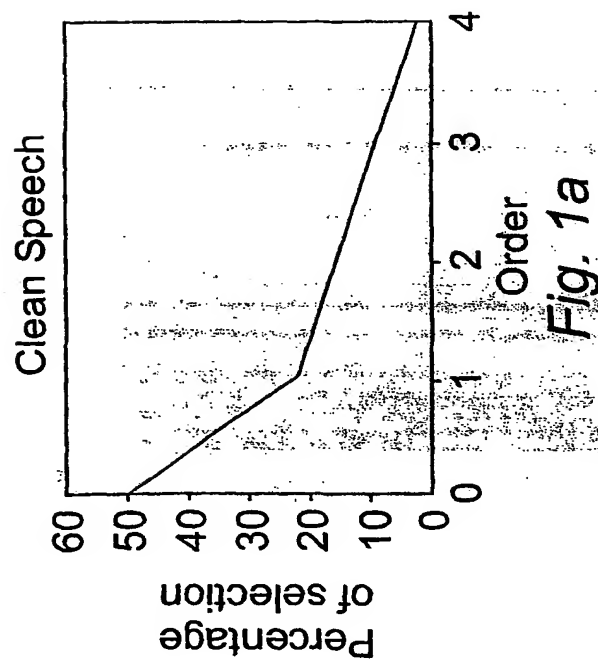
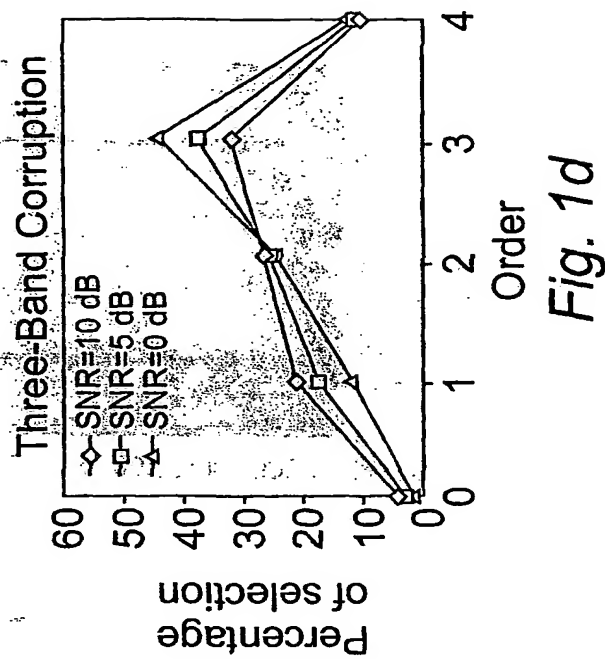
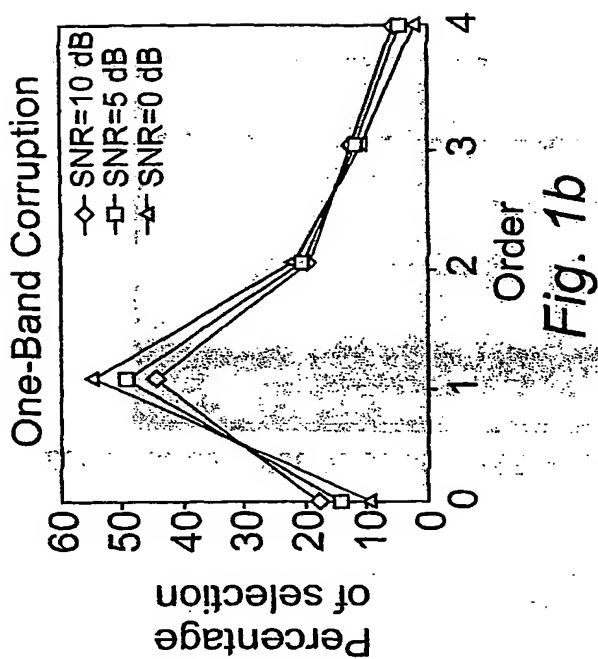
$U(r)$ is a recognition result;

$S(r)$ stands for the total number of states in $U(r)$;

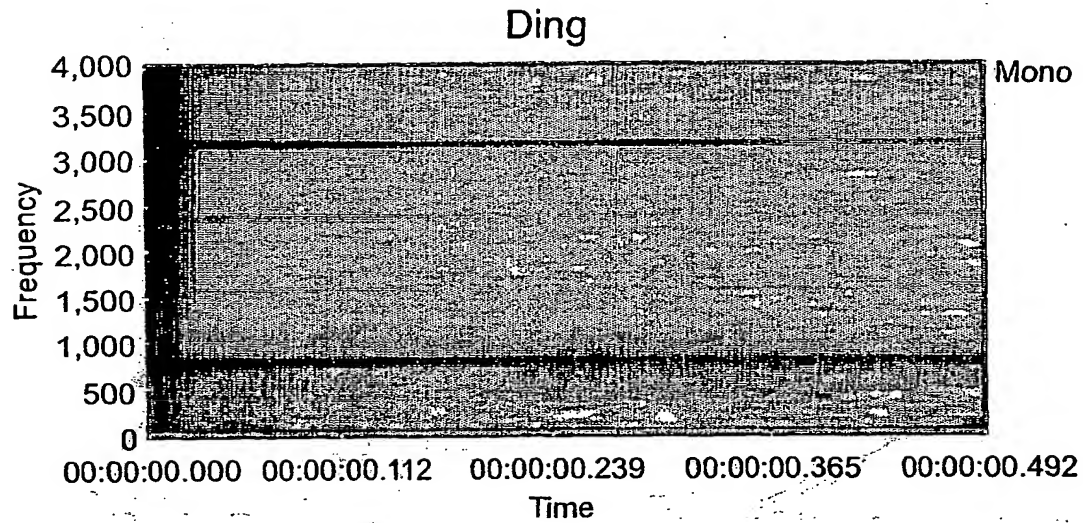
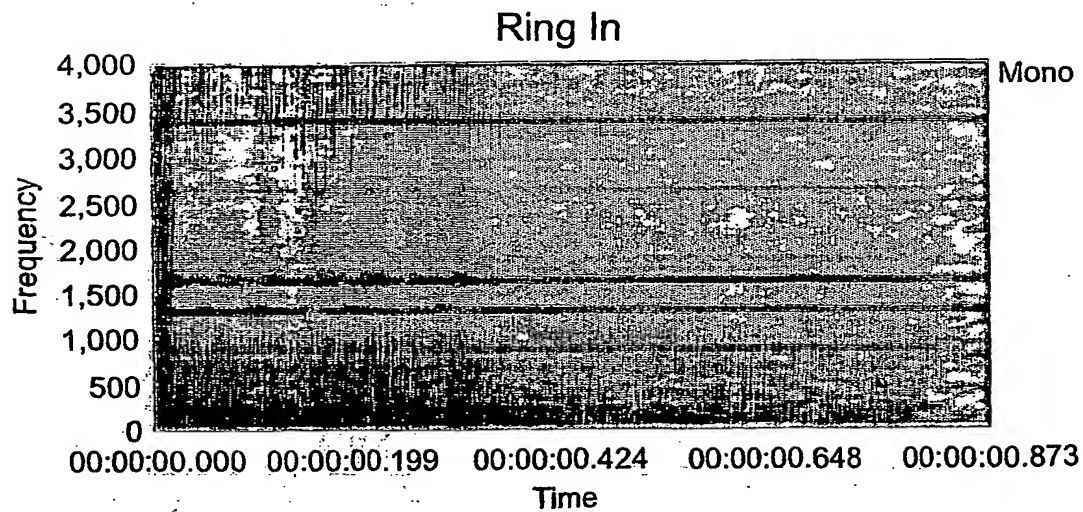
$P_i^u(d_i(r))$ is the state duration probability for d frames in state i of phonetic unit u .

8. A method as claimed in any preceding Claim in combination with conventional signal filtering techniques which remove known stationary corruptions.
9. A method as claimed in any of the preceding Claims substantially as hereinbefore described with reference to the accompanying tables and drawings.

1 / 4



2 / 4

*Fig. 2a**Fig. 2b*

SUBSTITUTE SHEET (RULE 26)

3 / 4

Whistle

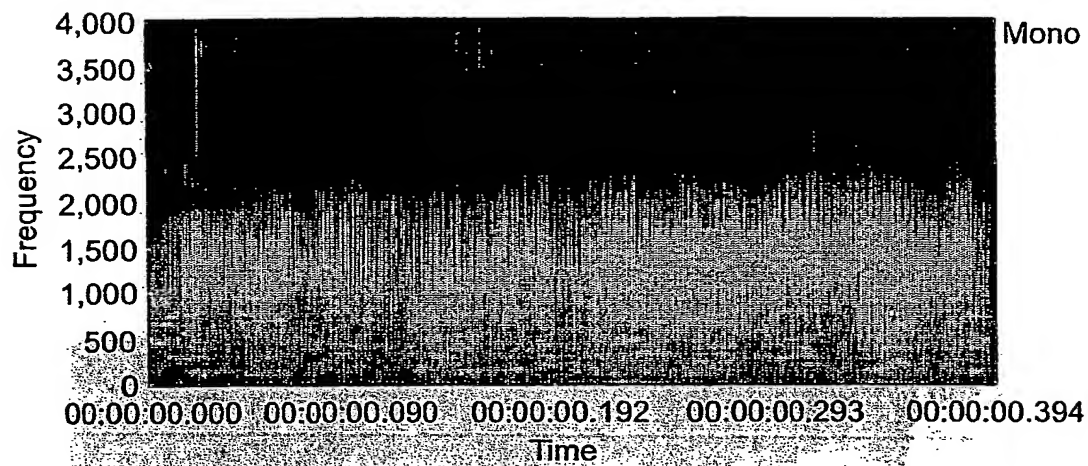


Fig. 2c

Contact

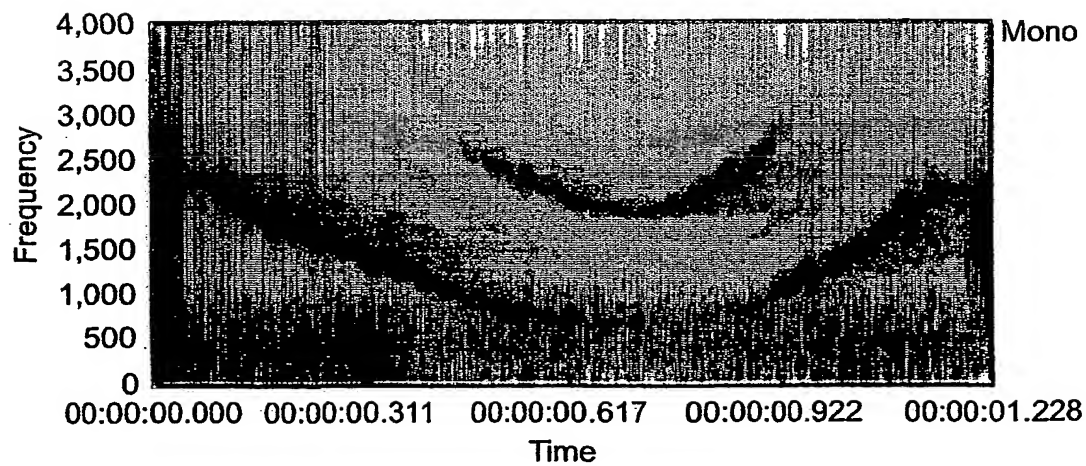


Fig. 2d

4 / 4

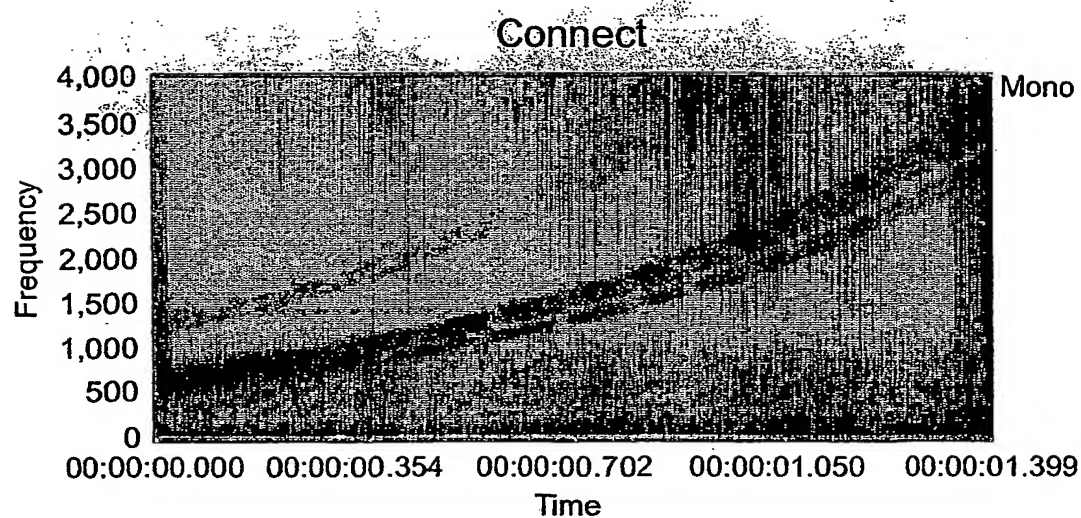


Fig. 2e

INTERNATIONAL SEARCH REPORT

PCT/GB 02/02197

A. CLASSIFICATION OF SUBJECT MATTER
IPC 7 G10L15/20

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)
IPC 7 G10L

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

EPO-Internal, WPI Data, INSPEC

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
P, X	JANCOVIC P ET AL: "A probabilistic union model with automatic order selection for noisy speech recognition" JOURNAL OF THE ACOUSTICAL SOCIETY OF AMERICA, SEPT. 2001, ACOUST. SOC. AMERICA THROUGH AIP, USA, vol. 110, no. 3, pages 1641-1648, XP001100608 ISSN: 0001-4966 the whole document ----- -/-	1-8

☒ Further documents are listed in the continuation of box C.

☐ Patent family members are listed in annex.

* Special categories of cited documents :

- *A* document defining the general state of the art which is not considered to be of particular relevance
- *E* earlier document but published on or after the international filing date
- *L* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)
- *O* document referring to an oral disclosure, use, exhibition or other means
- *P* document published prior to the international filing date but later than the priority date claimed

- *T* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
- *X* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
- *Y* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.
- *G* document member of the same patent family

Date of the actual completion of the international search

9 August 2002

Date of mailing of the international search report

22/08/2002

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax (+31-70) 340-3016

Authorized officer

Quélavoine, R

INTERNATIONAL SEARCH REPORT

PCT/GB 02/02197

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	JI MING ET AL: "Union: a new approach for combining sub-band observations for noisy speech recognition" SPEECH COMMUNICATION, APRIL 2001, ELSEVIER, NETHERLANDS, vol. 34, no. 1-2, pages 41-55, XP002209287 ISSN: 0167-6393 the whole document	1-5,8,9
X	JI MING ET AL: "A probabilistic union model for sub-band based robust speech recognition" 2000 IEEE INTERNATIONAL CONFERENCE ON ACOUSTICS, SPEECH, AND SIGNAL PROCESSING. PROCEEDINGS (CAT. NO.00CH37100), 5 - 9 June 2000, pages 1787-1790 vol.3, XP002209288 ISTANBUL, TURKEY, Piscataway, NJ, USA, IEEE, USA ISBN: 0-7803-6293-4 the whole document	1-5,8,9
A	JANCOVIC P ET AL: "Combining multi-band and frequency-filtering techniques for speech recognition in noisy environments" TEXT, SPEECH AND DIALOGUE. THIRD INTERNATIONAL WORKSHOP, TSD 2000. PROCEEDINGS (LECTURE NOTES IN ARTIFICIAL INTELLIGENCE VOL.1902), 13 - 16 September 2000, pages 265-270, XP008006658 BRNO, CZECH REPUBLIC, Berlin, Germany, Springer-Verlag, Germany ISBN: 3-540-41042-2 the whole document	1-9

Form PCT/ISA/210 (continuation of second sheet) (July 1992)